# Semi-supervised Sequential Generative Models

**Michael Teng**[1*], **Tuan Anh Le**[2], **Adam Scibior**[3], **Frank Wood**[3,4]

[1] Department of Engineering Science, University of Oxford
[2] Department of Brain and Cognitive Sciences, MIT
[3] Department of Computer Science, University of British Columbia
[4] Montréal Institute for Learning Algorithms (MILA)

## Abstract

We introduce a novel objective for training deep generative time-series models with discrete latent variables for which supervision is only sparsely available. This instance of semi-supervised learning is challenging for existing methods, because the exponential number of possible discrete latent configurations results in high variance gradient estimators. We first overcome this problem by extending the standard semi-supervised generative modeling objective with reweighted wake-sleep. However, we find that this approach still suffers when the frequency of available labels varies between training sequences. Finally, we introduce a unified objective inspired by teacher-forcing and show that this approach is robust to variable length supervision. We call the resulting method caffeinated wake-sleep (CWS) to emphasize its additional dependence on real data. We demonstrate its effectiveness with experiments on MNIST, handwriting, and fruit fly trajectory data.

## 1 INTRODUCTION

In recent years there has been an explosion of interest in deep generative models (DGM), which use neural networks transforming random inputs to learn complex probability distributions. We particularly focus on the variational auto-encoder (VAE ; Kingma and Welling (2013)) family of models, where the generative model is learned simultaneously with an associated inference network. This approach has been extended to settings with partially observed discrete variables (Kingma et al., 2014) and sequential data (Chung et al., 2015) but so far little work has been done on combining the two. In this paper we address this gap.

In many scenarios it is natural to assign discrete labels to specific intervals of time-series data and try to infer these labels from observations. For example, sequences of stock prices can be identified as periods of bear or bull markets, fragments of home CCTV footage can be classified as burglaries or other events, and heart rate signals can be used to classify cardiac health. In addition to the latent classification of observations in these applications, we may also wish to conditionally generate new sequences or predict ahead given some input sequence. Existing approaches to generative modeling in this context either require full supervision in the label space, and therefore are not able to leverage large amounts of unlabeled data, or are restricted to classification and can not generate new data (Chen et al., 2013; Wei and Keogh, 2006).

Additionally, many methods for semi-supervised learning of static DGM tasks define separate unsupervised and supervised loss terms in the training objective, which are weighted to reflect the overall supervision rate in the dataset (Kingma et al., 2014). Though this class of approaches is extendable to the time-series setting, the optimization of model and inference network parameters become unstable when there is an uneven contribution of partial labels to the supervised term, as a result of varying supervision rates in each training example. Because of this, the state-of-the-art approaches to semi-supervised learning, including Virtual Adversarial Training (Miyato et al., 2018), Mean Teacher (Tarvainen and Valpola, 2017), and entropy minimization (Grandvalet and Bengio, 2005), cannot be easily adapted to modeling time-series data.

In this work we derive two novel objectives for training DGMs with discrete latent variables in a semi-supervised fashion. They are both based on the reweighted wake-sleep (RWS) algorithm (Bornschein and Bengio, 2014), which avoids using high variance score function estimators for expectations over the discrete latents. The first one, which we call semi-supervised wake-sleep (SSWS), is obtained by extending RWS with a supervised classifi-

---

*Correspondence to `mteng@robots.ox.ac.uk`

cation term. While it performs well when labels available in the training set are regularly distributed per training example, it suffers from optimization problems when they are not, which can be unavoidable in real world datasets. To overcome this problem, we introduce a new approach which performs wake-sleep style optimization using a single objective that incorporates both supervised and unsupervised terms for every gradient update. We call this approach caffeinated wake-sleep (CWS).

Although we are explicitly targeting the sequential setting, our objectives can also be useful in non-sequential settings, especially when the discrete latent space is too large to be fully enumerated. We evaluate our objectives on MNIST, a handwriting dataset, and a dataset of fruit fly trajectories labeled with behavior classes. In the rest of the paper, Section 2 reviews the relevant background information, Section 3 describes our model and training objective, and Section 4 presents the experimental results.

## 2 BACKGROUND

### 2.1 VARIATIONAL AUTO-ENCODERS

Variational auto-encoders (Kingma and Welling, 2013; Rezende et al., 2014) are an approach to deep generative modeling where we simultaneously learn a generative model $p_\theta$ and an amortized inference network $q_\phi$, both parameterized by neural networks. For any given observation $\mathbf{x}$, the inference network variationally approximates the intractable posterior distribution $p_\theta(\mathbf{z}|\mathbf{x})$ over a latent variable $\mathbf{z}$ as $q_\phi(\mathbf{z}|\mathbf{x})$. In order to train both networks simultaneously, Kingma and Welling (2013) propose maximizing the evidence lower bound (ELBO), which is a sum over ELBOs for individual data points defined as

$$\mathcal{L}(\theta, \phi, \mathbf{x}) := \log p_\theta(\mathbf{x}) - \mathrm{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p_\theta(\mathbf{z}|\mathbf{x}))$$
$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \quad (1)$$

The ELBO approximates log-marginal likelihood $\log p_\theta(\mathbf{x})$, making it a good target objective for learning $\theta$, and is proportional to negative $\mathrm{KL}(q_\phi|p_\theta)$, making it a good target objective for learning $\phi$.

Burda et al. (2015) propose an extension to the variational autoencoder (VAE) where, for a given number of samples $K$, the single-datapoint objective is instead defined as

$$\mathcal{L}^K(\theta, \phi, \mathbf{x}) = \mathbb{E}_{\mathbf{z}_{1:K} \sim q_\phi} \left[ \log \left( \frac{1}{K} \sum_{k=1}^K \frac{p_\theta(\mathbf{z}_k, \mathbf{x})}{q_\phi(\mathbf{z}_k|\mathbf{x})} \right) \right] \quad (2)$$

This yields a tighter bound for $\log p_\theta(\mathbf{x})$, which is desirable for learning $\theta$.

In either case, the estimator of the gradient with respect to $\phi$ is typically high variance if $\mathbf{z}$ includes any discrete variables and therefore not suitable for gradient-based optimization. Various authors have proposed to reduce this variance using continuous relaxation Maddison et al. (2016); Jang et al. (2016) or control-variate methods (Mnih and Rezende, 2016; Mnih and Gregor, 2014; Tucker et al., 2017; Grathwohl et al., 2017) with varying degrees of success.

An alternative approach is to use the reweighted wake-sleep (RWS) algorithm (Bornschein and Bengio, 2014) which uses separate objectives for $\theta$ and $\phi$, interleaving the corresponding gradient steps. The target for $\theta$ is $\mathcal{L}_{\mathrm{IWAE}}^K(\theta, \phi, \mathbf{x})$, while the target for $\phi$ is $-\mathrm{KL}(p_\theta \parallel q_\phi)$. In the latter the expectation with respect to $p_\theta$ is approximated by importance sampling from $q_\phi$. Le et al. (2018) show that using two separate objectives avoids the problems with high variance gradient estimates in the presence of discrete latent variables and avoids learning problems discussed by Rainforth et al. (2018). Ba et al. (2015) demonstrate RWS as a viable method for time-series modeling by using it to train recurrent attention models.

### 2.2 SEMI-SUPERVISED VAE

In certain situations it is desirable to extend the VAE with an interpretable latent variable $y$, the canonical example being learning on the MNIST dataset where $y$ is the digit label, $\mathbf{z}$ is "style," and $\mathbf{x}$ is the image. Kingma et al. (2014) consider the setting where the label $y$ is only scarcely available and the dataset consists of many instances of $\mathbf{x}$ and relatively few instances of $(\mathbf{x}, y)$, which is a typical semi-supervised learning setting. They choose a factorization $p_\theta(\mathbf{x}, y, \mathbf{z}) = p_\theta(\mathbf{x}|y, \mathbf{z})p_\theta(y)p_\theta(\mathbf{z})$ and $q_\phi(y, \mathbf{z}|\mathbf{x}) = q_\phi(\mathbf{z}|\mathbf{x}, y)q_\phi(y|\mathbf{x})$. A naive optimization objective can be constructed by summing standard ELBOs for $\mathbf{x}$ and $(\mathbf{x}, y)$. However, this objective is unsatisfactory, since it does not contain any supervised learning signal for $q_\phi(y|\mathbf{x})$. Letting $\mathbf{x}_u$ denote unsupervised input variables and $(\mathbf{x}_s, y_s)$ denote supervised observation-latent pairs, Kingma et al. (2014) propose to maximize the corresponding supervised and unsupervised objectives:

$$\mathcal{L}^{\mathrm{s}}(\mathbf{x}_s, y_s) = \mathbb{E}_{q_\phi(\mathbf{z}|y_s, \mathbf{x}_s)} \left[ \log \left( \frac{p(\mathbf{x}_s, y_s, \mathbf{z})}{q_\phi(\mathbf{z}|y_s, \mathbf{x}_s)} \right) \right]$$
$$+ \alpha \, \mathbb{E}_{\hat{p}(y_s, \mathbf{x}_s)} \left[ \log(q_\phi(y_s|\mathbf{x}_s)) \right] \quad (3)$$

$$\mathcal{L}^{\mathrm{u}}(\mathbf{x}_u) = \mathbb{E}_{q_\phi(\mathbf{z}, y|\mathbf{x}_u)} \left[ \log \left( \frac{p(\mathbf{x}_u, y, \mathbf{z})}{q_\phi(\mathbf{z}, y|\mathbf{x}_u)} \right) \right] \quad (4)$$

The additional second term in Equation 3 is an expectation over the empirical distribution of labeled training pairs,
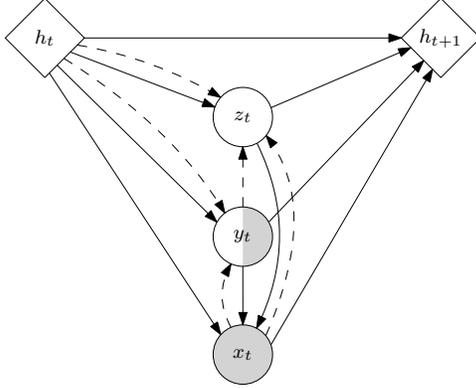
Figure 1: Per time step graphical model for the structured VRNN we use in this paper. In addition to the continuous latent variable $\mathbf{z}_t$, we include a discrete latent variable $y_t$ corresponding to an interpretable label, for which ground truth is sometimes available in the dataset. Dashed lines indicate the inference inference network, and solid lines the generative model.

$\hat{p}(\mathbf{x}_s, y_s)$, and can be optionally scaled by a hyperparameter, $\alpha$, controlling the relative strength of supervised and unsupervised learning signals. A semi-supervised VAE typically achieves performance superior to its unsupervised version. We subsequently refer to this objective as "M1+M2".

## 2.3  VARIATIONAL RNN

The variational RNN (VRNN) (Chung et al., 2015) is a deep latent generative model that is an extension of the VAE family. It can be viewed as an instantiation of a VAE at each time-step, with the model factorised in time overall. We use a structured VRNN variant, where a discrete latent variable $y$ is also present at each time step, as depicted in Figure 1.

## 2.4  RELATED WORK

Xu et al. (2017) introduce a semi-supervised VAE which overcomes the high-variance gradient estimators in sequential models using control variates during training. We build upon their work by adapting a similar architecture specification for the inference networks while avoiding the high-variance gradient estimators using reweighted wake-sleep style updates for training. Additionally, Chen et al. (2018) proposes a semi-supervised DGM for modeling natural language using full sentence context. While a useful method for the task domain, they state their model cannot do generation.

## 3  CAFFEINATED WAKE-SLEEP

Here we derive a general purpose training objective used for semi-supervised generative modeling. As we will show in subsequent experiments, both the objective and training method can be used to effectively train sequential models. As such, we introduce CWS in the context of time-series modeling, but note that it can trivially be used for static models by considering them to have a single time-step. In all subsequent sections, we denote $a_{\leq t} := a_{1:t}$ and $b_{<t} := b_{1:t-1}$ if $t > 1$.

### 3.1  GENERATIVE MODEL

First, we define a generative model over a sequence of observed variables, $\mathbf{x}_{1:T}$, continuous latent variables, $\mathbf{z}_{1:T}$, and a discrete latent categorical variable, $y_{1:T}$.

$$p_\theta(\mathbf{x}_{\leq T}, y_{\leq T}, \mathbf{z}_{\leq T}) = \prod_{t \leq T} \Big[ p_\theta(\mathbf{x}_t | \mathbf{z}_{\leq t}, y_{\leq t}) \qquad (5)$$
$$\times\, p_\theta(\mathbf{z}_t | \mathbf{x}_{<t}, y_{<t}, \mathbf{z}_{<t})\, p_\theta(y_t | \mathbf{x}_{<t}, y_{<t}, \mathbf{z}_{<t}) \Big]$$

Additionally, we are given $\mathbf{x}_{\leq T}$ and a subset $S \subseteq 1 : T$ of labeled $y_t$, denoted $y_S := \{y_t : t \in S\}$. We denote unlabeled $y_t$ as $y_U := \{y_t : t \in (1 : T) \setminus S\}$.

### 3.2  INFERENCE

Given the generative model defined above and a partially labeled dataset in the $y$ space, our variational distribution is the following:

$$q_\phi(y_U, \mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_S) = \qquad (6)$$
$$\prod_{t \in U} q_\phi(y_t | \mathbf{x}_{\leq T}, y_{<t}, \mathbf{z}_{<t}) \prod_{t \leq T} q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq T}, y_{\leq t}, \mathbf{z}_{<t})$$

If the categorical variable $y_t$ is always treated as a latent variable in the fully unsupervised case, we can define a variational distribution, $q(\mathbf{z}_{\leq T}, y_{\leq T} | \mathbf{x}_{\leq T})$ and maximize the ELBO with respect to $\{\theta, \phi\}$. Instead, we will not need to infer given $y_S$ for any sequence, but at any time-step when $y_t$ is available, we use it instead of a sampled $y_t$ from the variational distribution. This is made clear by writing the ELBO as:

$$\log p(\mathbf{x}_{\leq T}, y_S) \geq \qquad (7)$$
$$\mathbb{E}_{q_\phi(y_U, \mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_S)} \Bigg[ \log \Bigg( \prod_{t \in S} \frac{p_\theta(\mathbf{x}_t, \mathbf{z}_t, y_t,)}{q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq T}, y_{\leq t}, \mathbf{z}_{<t})}$$
$$\times \prod_{t \in U} \frac{p_\theta(\mathbf{x}_t, \mathbf{z}_t, y_t)}{q_\phi(y_t, \mathbf{z}_t | \mathbf{x}_{\leq T}, \mathbf{z}_{<t}, y_{<t})} \Bigg) \Bigg]$$

### 3.3  OPTIMIZATION

Having defined the model and inference network, we now specify the optimization objective.
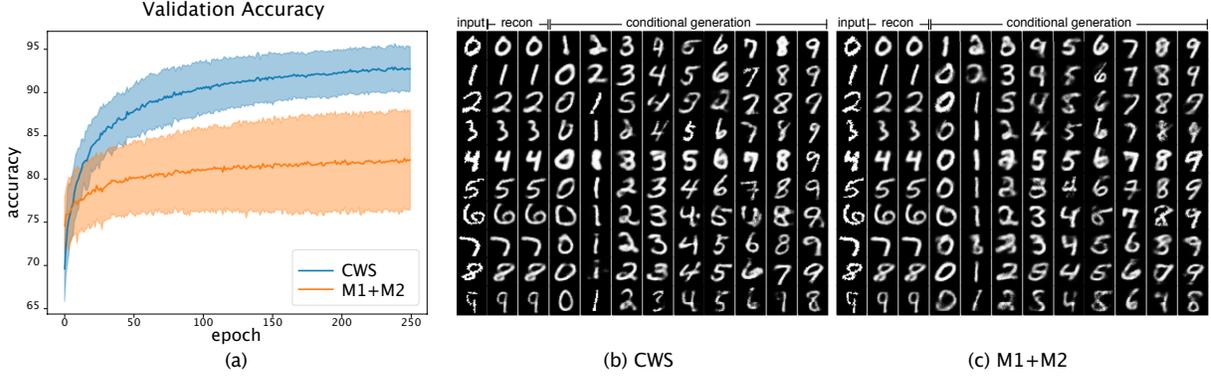
Figure 2: (a) Validation accuracy of a MNIST classification model comparing two approaches for semi-supervised learning, CWS and M1+M2. As seen, CWS trains to the highest validation accuracy and faster on a held out test set of 10000 examples. (b) Visual reconstructions of MNIST digit in final trained model using CWS (c) Reconstructions using the final trained model by M1+M2 method. In (b-c), the first column denotes the input image, the second denotes reconstruction where $z, y$ are both inferred using $q(y, z|x)$, the third denotes reconstruction where $z$ is inferred given the correct label $y$, and columns 4-12 denote reconstruction where we infer the style, $z$, fix it, and vary the $y$ label to conditionally generate.

### 3.3.1 SEMI-SUPERVISED WAKE-SLEEP

Before introducing our new objective, we first show how to use the reweighted wake-sleep algorithm (Bornschein and Bengio, 2014) with the semi-supervised objective introduced by Kingma et al. (2014). Following the recommendation of Le et al. (2018), extending reweighted wake-sleep to include a supervision term can be used to effectively avoid problems with discrete variables. Below we derive the semi-supervised wake-sleep (SSWS) objective within that framework. Both of our objectives will be evaluated by sampling from $q_\phi$, so for each $k \in 1 : K$ we have $y_t^k \sim q_\phi(y_t|\cdot)$ for all $t \in U$ and $\mathbf{z}_t^k \sim q_\phi(\mathbf{z}_t|\cdot)$ for all $t \in 1 : T$ with a convention $y_t^k = y_t$ for $t \in S$.

For learning the generative model, we maximize the IWAE bound with respect to parameters $\theta$:

$$\mathcal{L}_p(\mathbf{x}_{\leq T}, y_S) :=$$
$$\mathbb{E}_q \left[ \log \frac{1}{K} \sum_{k=1}^{K} \frac{p_\theta(y_U^k, \mathbf{z}_{\leq T}^k, \mathbf{x}_{\leq T}, y_S)}{q_\phi(y_U^k, \mathbf{z}_{\leq T}^k | \mathbf{x}_{\leq T}, y_S)} \right], \quad (8)$$
$$q = \prod_{k=1}^{K} q_\phi(y_U^k, \mathbf{z}_{\leq T}^k | \mathbf{x}_{\leq T}, y_S)$$

This is a lower bound to $\log p_\theta(\mathbf{x}_{\leq T}, y_S)$ which is tight when $q_\phi(y_U, \mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_S) = p_\theta(y_U, \mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_S)$. Sampling from $q_\phi(y_U, \mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_S)$ and evaluating its density is simple since it is factorized as in (6) where both $q(y_t|\cdot)$ and $q_\phi(\mathbf{z}_t|\cdot)$ are given, and all values to the right of the conditioning bar are always available at sampling at time step $t$ (previous $y_{<t}$ are either sampled or given supervision and $\mathbf{z}_{<t}$ are sampled).

For learning the inference network parameters, we need to consider the unsupervised and the supervised cases. For unsupervised inference learning, we minimize the expected KL-divergence between the true posterior and the variational posterior given by $q_\phi$ under the generative model, $p := p_\theta(\mathbf{x}_{\leq T}, y_S)$:

$$\nabla_\phi \mathbb{E}_p[\text{KL}(p_\theta(y_U, \mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_S) || q_\phi(y_U, \mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_S))]$$
$$= \mathbb{E}_p[\mathbb{E}_{p_\theta(y_U, \mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_S)} [-\nabla_\phi \log q_\phi(y_U, \mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_S)]]$$

Given some $(\mathbf{x}_{\leq T}, y_S)$ from the data distribution, $\hat{p}_\theta(\mathbf{x}_{\leq T}, y_S)$, the inner expectation is approximated using samples from the inference network, $q_\phi$, referred to as the wake-$\phi$ update for learning parameters $\phi$:

$$\mathcal{L}_q(\mathbf{x}_{\leq T}, y_S) :=$$
$$\mathbb{E}_{q_\phi} \left[ \sum_{k=1}^{K} \bar{w}_k \left( -\log q_\phi(y_U^k, \mathbf{z}_{\leq T}^k | \mathbf{x}_{\leq T}, y_S) \right) \right], \quad (9)$$

where $\bar{w}_k$ is a normalized version of

$$w_k := \frac{p_\theta(y_U^k, \mathbf{z}_{\leq T}^k, \mathbf{x}_{\leq T}, y_S)}{q_\phi(y_U^k, \mathbf{z}_{\leq T}^k | \mathbf{x}_{\leq T}, y_S)}, \quad \bar{w}_k = \frac{w_k}{\sum_l w_l} \quad (10)$$

and

$$q_\phi = \prod_k q_\phi(y_U^k, \mathbf{z}_{\leq T}^k | \mathbf{x}_{\leq T}, y_S)$$

Evaluating (9) requires evaluating the joint $p$ density given in (5), the $q$ density given in (6), and sampling again from (6).

Finally, we need to include a supervised loss term. There is only supervision signal for some given $y_t$ for some $t \in S$ per data example. As such, we want to maximize
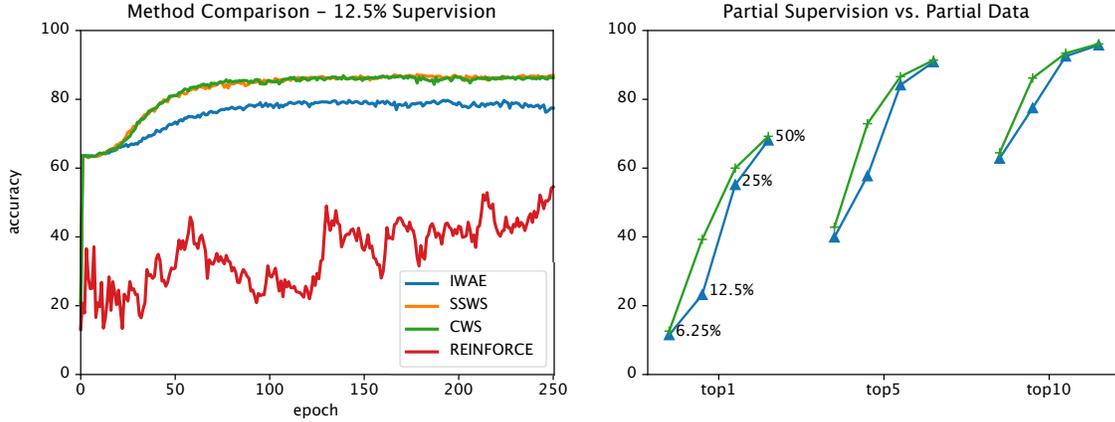
Figure 3: Left: Top10 validation accuracy throughout training for a 70-class character classification task comparing three methods of semi-supervised training, CWS, SSWS, and REINFORCE (M1+M2), with 12.5% supervision rate and one method of fully supervised training using IWAE and a classifier loss on 12.5% of the dataset. Right: Final top1, top5, and top10 accuracy summaries for semi-supervision using CWS (green) against models trained with full supervision but on a corresponding fraction of the original dataset (blue). As before, the method denoted by blue trains with full supervision using IWAE and a classifier loss on the subset of the total data denoted by the labels.

the log-likelihood of the supervised pairs $(y_S, \mathbf{x}_{\leq T})$ under the variational distribution, $q_\phi(y_{\leq T}|\mathbf{x}_{\leq T})$, marginalized over unsupervised time-steps:

$$q_\phi(y_S|\mathbf{x}_{\leq T}) = \int q_\phi(y_{\leq T}, \mathbf{z}_{\leq T}|\mathbf{x}_{\leq T}) dy_U d\mathbf{z}_{\leq T}$$

$$= \int q_\phi(y_S|\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}, y_U) q_\phi(y_U, \mathbf{z}_{\leq T}|\mathbf{x}_{\leq T}) dy_U d\mathbf{z}_{\leq T}$$

$$= \mathbb{E}_{q_\phi(y_U, \mathbf{z}_{\leq T}|\mathbf{x}_{\leq T})}[q_\phi(y_S|\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}, y_U)] \quad (11)$$

We lower bound this with Jensen's inequality and obtain our supervised loss term:

$$\mathcal{L}_s(y_S) :=$$
$$\mathbb{E}_{q_\phi(y_U, \mathbf{z}_{\leq T}|\mathbf{x}_{\leq T})}[\log q_\phi(y_S|\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}, y_U)] \quad (12)$$

For each minibatch the objectives, (8), (9), and (12), are computed using the same set of samples and we perform the relevant gradient steps in $\phi$ and $\theta$ by alternating between the two parameter updates:

$$\theta^* = \theta + \alpha_\theta \nabla_\theta \mathcal{L}_p(\mathbf{x}_{\leq T}, y_S) \quad (13)$$

$$\phi^* = \phi - \alpha_\phi \nabla_\phi \Big( \mathcal{L}_q(\mathbf{x}_{\leq T}, y_S) - \mathcal{L}_s(y_S) \Big) \quad (14)$$

### 3.3.2 CAFFEINATED WAKE-SLEEP

Although the SSWS approach can be used to train time-series models, in practice there is a tradeoff between learning and optimization stability. As mentioned earlier, sequences of observations may contain variable amounts of supervision, the most extreme example being datasets

containing both sequences with fully observed and fully unobserved labels.

Because of this, the magnitude of the supervised and unsupervised terms in the $q_\phi$ loss will vary per data sequence, and we incur a tradeoff between correcting for this bias and computational efficiency. For example, we could normalize $\mathcal{L}_s$ and $\mathcal{L}_q$ by $|S|$ and $|U|$, respectively, but doing so treats sequences with a lower supervision rate the same as those with much higher supervision. Alternatively, we could take a weighted average of the terms across the sequences, but this faces the same issue as before unless we also scale the learning rate dynamically per gradient step for each stochastic mini-batch.

To remedy this problem, we obviate the need to maintain two different supervision terms by minimizing the expected KL-divergence between the true posterior and the full variational posterior under the generative likelihood, $p_\theta(\mathbf{x}_{\leq T})$:

$$\nabla_\phi \mathbb{E}_{p_\theta(\mathbf{x}_{\leq T})}[\text{KL}(p_\theta(y_{\leq T}, \mathbf{z}_{\leq T}|\mathbf{x}_{\leq T})||q_\phi(y_{\leq T}, \mathbf{z}_{\leq T}|\mathbf{x}_{\leq T}))]$$
$$= \mathbb{E}_{p_\theta(\mathbf{x}_{\leq T})}[\mathbb{E}_{p_\theta(y_U, y_S, \mathbf{z}_{\leq T}|\mathbf{x}_{\leq T})}[-\nabla_\phi \log q_\phi(\cdot|\mathbf{x}_{\leq T})]]$$

Unlike in the derivation for Equation (9), there is no distinction here between $y_S$ and $y_U$, meaning that we should always sample a value for $y_t$ when computing this loss. However, doing so would give the fully unsupervised wake-$\phi$ objective and not include any supervision signal for the $y_S$ labels that we do have. To "correct" for this, we introduce an empirically justified bias into this estimator, namely replacing sampled values of $y_t$ with $y_S$ when available. Intuitively, this is exactly the reweighted wake-
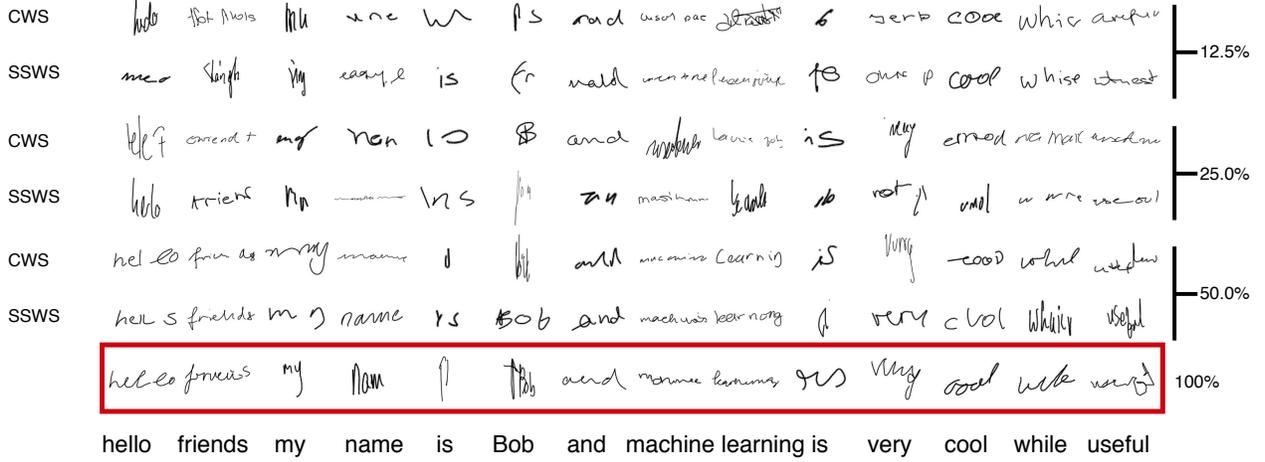
Figure 4: Generated of sample trajectories of handwriting conditioned on generating the bottom text. Each pair of rows is trained with a supervision rate (right) and a method (left). Corresponding sample generated from a fully supervised model is displayed in the red box. Note that the 100% supervision example in this figure is trained with the full dataset with a separate classifier loss.

sleep algorithm if $y_S$ is treated as having been sampled from the distributions, $q_\phi(\cdot)$. This is a kind of teacher-forcing approach (Williams and Zipser, 1989) for discrete latent labels, $y_t$.

To be concrete about our method, we still compute an expectation over the empirical data distribution $\hat{p}(\mathbf{x}_{\leq T}, y_S)$, and estimate the innermost expectation with samples from (6) when needed. The bias introduced then comes from the additional $q_\phi(y_S|...)$ terms in the denominator of the importance weights:

$$\mathcal{L}_q^{\mathrm{CWS}}(\mathbf{x}_{\leq T}, y_S) :=$$

$$\mathbb{E}_{(y_S, \mathbf{x}_{\leq T}) \sim \hat{p}, \ (y_U^{1:K}, \mathbf{z}^{1:K}) \sim q_\phi(y_U^k, \mathbf{z}_{\leq T}^k | y_S, \mathbf{x}_{\leq T})} [f],$$

$$f := \sum_{k=1}^K \tilde{w}_k \left( -\log q_\phi(y_U^k, y_S, \mathbf{z}_{\leq T}^k | \mathbf{x}_{\leq T}) \right), \qquad (15)$$

$$w_k^{\mathrm{cws}} := \frac{p_\theta(y_U^k, \mathbf{z}_{\leq T}^k, \mathbf{x}_{\leq T}, y_S)}{q_\phi(y_U^k, y_S, \mathbf{z}_{\leq T}^k | \mathbf{x}_{\leq T})}, \quad \tilde{w}_k = \frac{w_k^{\mathrm{cws}}}{\sum_l w_l^{\mathrm{cws}}}$$

Finally, we formally introduce the caffeinated wake-sleep algorithm. Model parameters are learned with the importance-weighted ELBO as before. However, crucially, we use a unified gradient estimator in (15) to update $\phi$ parameters:

$$\theta^* = \theta + \alpha_\theta \nabla_\theta \mathcal{L}_p(\mathbf{x}_{\leq T}, y_S) \qquad (16)$$

$$\phi^* = \phi - \alpha_\phi \nabla_\phi \mathcal{L}_q^{\mathrm{CWS}}(\mathbf{x}_{\leq T}, y_S) \qquad (17)$$

## 4 EXPERIMENTS

We evaluate CWS on an variety of tasks within generative modeling by running experiments on three separate datasets, MNIST, IAM On-Line Handwriting Database (IAM-OnDB) (Liwicki and Bunke, 2005), and Fly-vs-Fly (Eyjolfsdottir et al., 2014). In our experiments, we consider the training accuracy of the inference network in classification, the conditional and unconditional generation of new data, and the uncertainty captured by our generative models.

### 4.1 SEMI-SUPERVISED MNIST

We start with a toy MNIST experiment, which semi-supervises a discrete latent variable corresponding to digit class. We use the same model and VAE architecture detailed in Kingma et al. (2014) and compare against their training method, denoted as M1+M2 with the supervised weight set to 60. For supervision, we use only 100 labeled digits out of a total dataset of 50000. We briefly note that CWS for the single time-step case is trivial given the objectives defined above and SSWS recovers the M1+M2 objective.

In Figure 2a, we compare validation accuracy of digit classification between CWS and M1+M2 trained models. We find that our method both trains faster and better despite using an identical generative model, inference network, learning rate, and optimizer. Furthermore, our approach requires setting one fewer hyperparameter when compared to M1+M2.

In Figure 2b, we conditionally generate samples using our trained models by first inferring the style using $q_\phi(\mathbf{z}|\hat{\mathbf{x}}, \hat{y})$. Then, fixing the style, we change $y$ and reconstruct using $p_\phi(\mathbf{x}|\mathbf{z}, y)$. We find that CWS is able to conditionally generate as well as prior art.
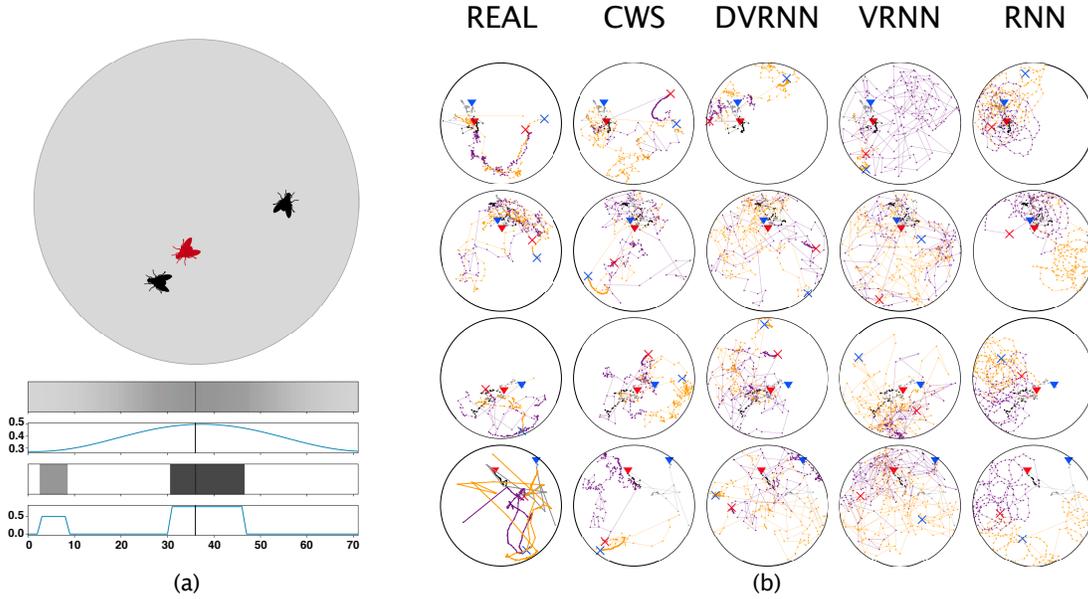
Figure 5: (a) Diagram of the red fly's field of view encoding in a petri dish environment containing two other flies. The middle line in the plots denotes the fly's direct line of sight. Each pair of plots indicates the agent's field of view with respect to walls and other flies, respectively. The fly closer to the red agent contributes more mass to the encoding vector. (b) Comparison of each model's continuation tracks for four distinct seed sequences. Each column shows a model's continuation sampled from the generative model compared against the leftmost column displaying the ground truth continuation. Each row shows one of four seed sequences used to seed the generative models. Within the petri dish, both flies' locomotion tracks are shown. For the first fly, the red arrow indicates the starting position of the seed sequence, the black markers and line indicate the seed tracks, the purple markers and line indicate the sampled continuation tracks, and the red 'x' indicates the final position of the fly after 200 time-steps. For the second fly, these indicators are blue arrow, gray markers, orange markers, and blue 'x', respectively.

## 4.2 HANDWRITING

Next, we run an experiment on the IAM-OnDB dataset. The generative model is over $\mathbf{x}_t := \{c_t^i, c_t^j, \mathrm{pen}_t, \mathrm{eoc}_t\}$, where $c_i, c_j$ denotes a single time-step vectorized stroke of the pen and the $\mathrm{pen}_t$ and $\mathrm{eoc}_t$ denote pen-up status and end-of-character binary values. The latent space in this model is comprised of continuous style variable, $\mathbf{z}_t$, and a discrete character label corresponding to a sequence of observed strokes that make up a valid character, $y_t$. For the IAM-OnDB data, the valid alphabet is comprised of 70 different characters: upper and lower case letters, digits, and special characters. The large latent space makes this a very challenging time-series problem.

The VRNN architecture we use is taken from the Deep-Writing model introduced by Aksan et al. (2018). At a high level, we use a VRNN over $\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}$, and a BiL-STM network for $q_\phi(y_{\leq T}|\mathbf{x}_{\leq T})$. Because the VRNN can be optimized solely using reparameterization, we train $q_\phi(\mathbf{z}_{\leq T}|y_{\leq T}, \mathbf{x}_{\leq T})$ using IWAE and $q_\phi(y_{\leq T}|\mathbf{x}_{\leq T})$ using CWS. This mirrors the optimization in Aksan et al. (2018), except that now, we can train with semi-supervision. For more detailed experimental setup, we refer the reader to

the cited work. In all experiments, we use a training set of 26560 sequences and a validation set of 512 sequences, all of length 200.

As previously mentioned, techniques for semi-supervised learning do not trivially extend to time-series models. Instead, we compare CWS against SSWS, the latter of which can be viewed as the M1+M2 objective for sequential models using wake-$\phi$ style updates for learning discrete latents. For comparison, we can also train the M1+M2 objective without using RWS by using the RE-INFORCE estimator for taking gradients through discrete latent variables. In Figure 3 (left), we find that this method results in poorly trained inference networks which cannot classify character labels accurately even using top-10 metric.

In Figure 3 (right), we present classification results of training with a fixed architecture, varying only the training objective using either CWS or IWAE. Because the normal IWAE objective cannot be used with semi-supervision and REINFORCE, we compare CWS with fully unsupervised IWAE plus training the classifier separately with full supervision. This baseline is generated by using the
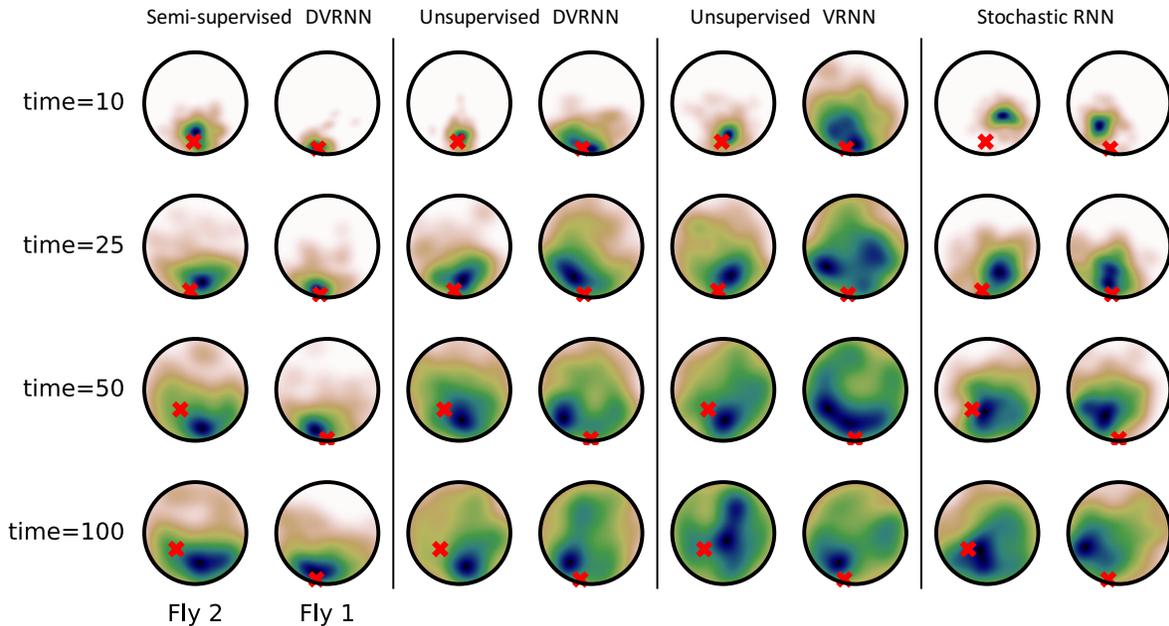
Figure 6: Comparison of uncertainty estimates over future trajectories of two flies' interactions in the environment between different models. Each model is seeded with the same ground truth sequence each time and 100 continuations are sampled for 100 time-steps into the future. We show the kernel density estimate of all fly positions at the indicated time-step across the 100 continuations. Each model's estimate is split by fly for clarity (two columns per model indicated by separators), where the red 'x' indicates the true position of the fly at that time-step. The future trajectories across all models show greater uncertainty about the fly's position as the model evolves in time. We find the model trained with CWS is the least noisy as it evolves, still captures the actual position within a high probability region, and has the most mass close to the true position of fly 1 at time 100 among the 4 models.

optimization technique used from Aksan et al. (2018) on a corresponding fraction of the dataset. In other words, we take the exact same architecture and compare training with conventional DGM techniques plus a classification loss against training with our CWS method using the same amount of full labels but being able to incorporate unlabeled data. We find that validation accuracy of the classifier trained with CWS using additional unlabeled data greatly outperforms the IWAE baseline using the same amount of partial labels from the dataset but without the additional data.

In Figure 4, we show reconstructed samples from trained models where we force the network to attempt to generate the phrase, "hello friends my name is Bob and machine learning is very cool while useful". We find that even at 12.5% supervision, sampled generations resemble the target sentence using SSWS or CWS.

### 4.3  FLY TRACKING

In our final experiment, we run CWS on a dataset of two male fruit flies interacting in a petri dish (Eyjolfsdottir

et al., 2014). We use a VRNN which attempts to mirror biological plausibility by modeling movement, $\mathbf{x}_t$, as velocities of body position and configuration (Cueva and Wei, 2018). We inject as input to the RNN at each time-step, the visual encoding for a given fly shown in Figure 5a (Clandinin and Giocomo, 2015) along with knowledge of its own state (Cueva and Wei, 2018).

The dataset uses high level actions, which we model as a 6-valued discrete variable, $y_t$, corresponding to possible semantic labels: "lunge", "charge", "tussle", "wing threat", "hold", and "unknown". These annotations are provided by human experts at time-steps dispersed throughout the dataset without a clear pattern or regularity. In all experiments we compare four models: the RNN used by Eyjolfsdottir et al. (2016) outputting a probability vector over discretized action space at each time step, the standard VRNN without $y$ and with continuous valued $\mathbf{z}$, the discrete VRNN (DVRNN) with $y_t$ trained unsupervised, and the same DVRNN trained with CWS. While CWS was able to train this model, we ran into optimization difficulties with the same model but training with a separate supervision term. We fine-tune and report results of the

Table 1: KDE of ground truth position under model

| | KDE $\log p$ | |
| model | fly 1 | fly 2 |
| --- | --- | --- |
| CWS+DVRNN | $-780.1$ | $-851.5$ |
| RWS+DVRNN | $-901.8$ | $-853.9$ |
| VRNN | $-962.0$ | $-879.8$ |
| RNN | $-931.0$ | $-921.2$ |

best performing model during training: CWS-DVRNN with 30-dimensional $z$ and 6-valued $y$, RWS-DVRNN with 50-dimensional $z$ and 10-valued $y$, and VRNN with 120-dimension $z$. We provide further details of model in the Appendix.

In Figure 5b, we condition the model on an initial sequence of actions, then sample a continuation and visually inspect how it compares with the true continuation that was not shown to the model. We find that under the RNN model, the flies tend to move in circular patterns with relatively constant velocity. In contrast, real flies tend to alternate between fast and slow movements, changing their directions much more abruptly. All three variants of our VRNN model qualitatively recover this behavior, however the continuous VRNN also tends to behave too erratically. In the DVRNN models, we were not able to identify any other clear visual artifacts in the generated trajectories that disagrees with the real data, although DVRNN without semi-supervision appears to generate trajectories where flies move too much. In Figure 6, we investigate the quality of uncertainty estimates produced by various models. For this purpose we again seed the model with an initial sequence of actions, then observe how the probability mass over the flies' projected future positions evolves over time and compare it with the actual positions in the dataset. Figure 6 visualizes the results and Table 1 provides the log likelihood of each fly's true position under this density estimate. Again, we find that CWS trained DVRNN performs best, followed by the unsupervised DVRNN, the continuous VRNN, and the RNN.

## 4.4 DISCUSSION

We have introduced a new method for semi-supervised learning in deep generative time-series models and we have shown that it achieves better performance than unsupervised learning and fully supervised learning with a fraction of the data. Although CWS uses a biased gradient estimator, it has many advantages including ease of implementation, intuition as teacher-forcing RWS, and empirical validation. A formal, theoretical justification of why the biased estimator of CWS works remains for future work. Additionally, we are interested in extending CWS to a more general class of semi-supervised models, taking inspiration from the ideas of Siddharth et al. (2017).

## References

Emre Aksan, Fabrizio Pece, and Otmar Hilliges. Deepwriting: Making digital ink editable via deep generative modeling. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, page 205. ACM, 2018.

Jimmy Ba, Ruslan R Salakhutdinov, Roger B Grosse, and Brendan J Frey. Learning wake-sleep recurrent attention models. In *Advances in Neural Information Processing Systems*, pages 2593–2601, 2015.

Jörg Bornschein and Yoshua Bengio. Reweighted wakesleep. *arXiv preprint arXiv:1406.2751*, 2014.

Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.

Mingda Chen, Qingming Tang, Karen Livescu, and Kevin Gimpel. Variational sequential labelers for semi-supervised learning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 215–226, 2018.

Yanping Chen, Bing Hu, Eamonn Keogh, and Gustavo EAPA Batista. Dtw-d: time series semi-supervised learning from a single example. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 383–391. ACM, 2013.

Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. A recurrent latent variable model for sequential data. In *Advances in neural information processing systems*, pages 2980–2988, 2015.

Thomas R Clandinin and Lisa M Giocomo. Neuroscience: Internal compass puts flies in their place. *Nature*, 521 (7551):165, 2015.

Christopher J Cueva and Xue-Xin Wei. Emergence of grid-like representations by training recurrent neural networks to perform spatial localization. *arXiv preprint arXiv:1803.07770*, 2018.

Eyrun Eyjolfsdottir, Steve Branson, Xavier P Burgos-Artizzu, Eric D Hoopfer, Jonathan Schor, David J Anderson, and Pietro Perona. Detecting social actions of fruit flies. In *European Conference on Computer Vision*, pages 772–787. Springer, 2014.

Eyrun Eyjolfsdottir, Kristin Branson, Yisong Yue, and Pietro Perona. Learning recurrent representations for hierarchical behavior modeling. In *International Conference on Learning Representations*, 2016.

Yves Grandvalet and Yoshua Bengio. Semi-supervised learning by entropy minimization. In *Advances in neural information processing systems*, pages 529–536, 2005.

Will Grathwohl, Dami Choi, Yuhuai Wu, Geoff Roeder, and David Duvenaud. Backpropagation through the void: Optimizing control variates for black-box gradient estimation. *arXiv preprint arXiv:1711.00123*, 2017.

Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.

Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.

Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In *Advances in neural information processing systems*, pages 3581–3589, 2014.

Tuan Anh Le, Adam R Kosiorek, N Siddharth, Yee Whye Teh, and Frank Wood. Revisiting reweighted wake-sleep. *arXiv preprint arXiv:1805.10469*, 2018.

Marcus Liwicki and Horst Bunke. Iam-ondb-an on-line english sentence database acquired from handwritten text on a whiteboard. In *Eighth International Conference on Document Analysis and Recognition (ICDAR'05)*, pages 956–961. IEEE, 2005.

Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.

Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018.

Andriy Mnih and Karol Gregor. Neural variational inference and learning in belief networks. *arXiv preprint arXiv:1402.0030*, 2014.

Andriy Mnih and Danilo J Rezende. Variational inference for monte carlo objectives. *arXiv preprint arXiv:1602.06725*, 2016.

Tom Rainforth, Adam R. Kosiorek, Tuan Anh Le, Chris J. Maddison, Maximilian Igl, Frank Wood, and Yee Whye Teh. Tighter variational bounds are not necessarily better. In *International Conference on Machine Learning*, 2018.

Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.

N. Siddharth, T Brooks Paige, Jan-Willem Van de Meent, Alban Desmaison, Noah Goodman, Pushmeet Kohli, Frank Wood, and Philip Torr. Learning disentangled representations with semi-supervised deep generative models. In *Advances in Neural Information Processing Systems*, pages 5925–5935, 2017.

Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *Advances in neural information processing systems*, pages 1195–1204, 2017.

George Tucker, Andriy Mnih, Chris J Maddison, John Lawson, and Jascha Sohl-Dickstein. Rebar: Low-variance, unbiased gradient estimates for discrete latent variable models. In *Advances in Neural Information Processing Systems*, pages 2627–2636, 2017.

Li Wei and Eamonn Keogh. Semi-supervised time series classification. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 748–753. ACM, 2006.

Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989.

Weidi Xu, Haoze Sun, Chao Deng, and Ying Tan. Variational autoencoder for semi-supervised text classification. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.